(12)                    **EUROPEAN PATENT APPLICATION**
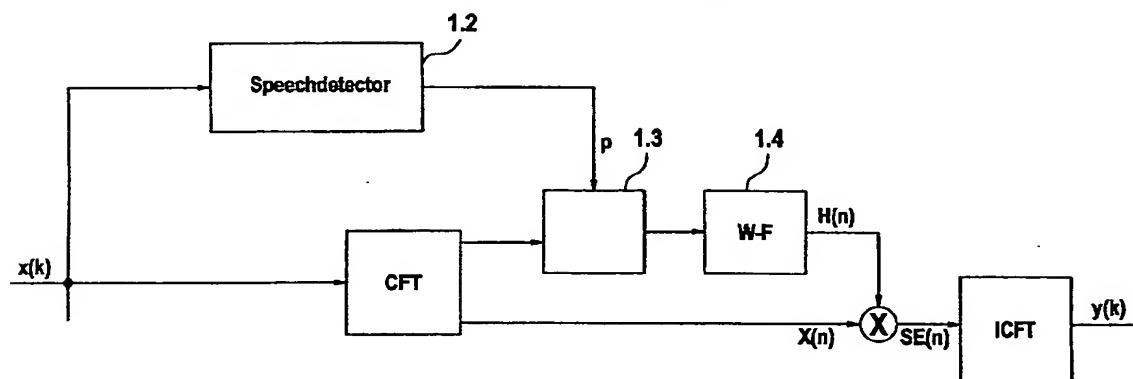
(54)    **Method for improving noise reduction in speech transmission**

(57)     Noise reduction measures must be taken in order to ensure a natural speech transmission in a noise-filled environment. This is particularly necessary in the case of speech-controlled appliances, in which speech recognition is an important quality feature. So-called spectral subtraction is used, as is known, for the purpose of noise reduction. In order to improve the determination of the noise components of a noisy speech signal using a Wiener filter, the conditions for calculation of the transmission function H(n) of the Wiener filter are adapted, according to the invention, to the nonlinear transmission behaviour of the human ear. For this purpose, in combination with the specified conditions, a Continuous Fourier Transformation is advantageously performed which prevents the occurrence of so-called musical tones. Despite a large noise reduction, loss of quality in the speech transmission is prevented by the method.

**Fig. 2**



EP 1 278 185 A2

## Description

[0001]   Where a speech signal is overlaid with unwanted noise it is essential to use methods for noise reduction. In the use of mobile telephones, unwanted noises are, for example, street noise, flight noise or noise in sports stadia. In order to ensure a natural speech transmission from a noise-filled environment, it is necessary to take measures to reduce the noise in the speech transmission. There is also an increasing use of speech-controlled appliances in which speech recognition is an important quality feature and which is essentially dependent on the mastery of noise reduction. The same problem must be resolved in the case of coding, for converting speech into text.

[0002]   DE 69 420 705 describes a system for noise suppression which comprises a multiplicity of microphones, signal processing means and an adaptive filter, which is preferably a Wiener filter. Auto and cross power spectra are determined from frequency-transformed sampling values of the speech signals. The signal processing means are provided in order to determine combined auto and cross power spectra from the auto and cross power spectra. The combined auto and cross power spectra provide the coefficients for the adaptive filter.

[0003]   DE 696 06 978 describes a method for noise suppression by means of spectral subtraction. In that case, non-speech frames are estimated using a non-parametric power spectrum estimation method, all N sampling values of each frame being used. A stationary background noise is assumed over several frames and a reduction of the variance of the power spectrum estimated value is achieved through averaging of the power spectrum estimated value over several non-speech frames. Speech frames are estimated using a parametric power spectrum estimation method, on the basis of a parametric model. Each speech frame contains a predefined number N of audio sampling values, as a result of which N degrees of freedom are assigned to each speech frame. The variance of the power spectrum estimation is reduced in that the parametric model contains few parameters, the parametric model reducing the number N of the degrees of freedom to the number of the parameters of the parametric model.

[0004]   A generally known method for noise reduction is that of so-called spectral subtraction. In this method, the noisy speech signal is first transformed from the time domain into the frequency domain, for example, by means of the Fast Fourier Transformation FFT, the noise spectrum is then determined in the speech pauses and subtracted from the frequency spectrum of the noisy speech signal before the noisy speech signal is reconverted from the from the frequency domain into the time domain by means of the Inverse Fast Fourier Transformation IFFT. The result depends essentially on the accuracy of the determination of the noise spectrum. Although good results are achieved in the case of stationary noise, in practice noises are not stationary and the achievable results are therefore unsatisfactory.

[0005]   Methods for spectral subtractions are described, for example, in the publications "Improved Spectral Subtraction for Speech Enhancement", Y. Malca, D. Wulich, and "Extended Spectral Subtraction", P. Sovka, P. Pollak, J. Kubie; EUSIPCO '96 Proceedings, Trieste, 10 - 13 September'96. These publications also make reference to fundamental works relating to spectral subtraction.

[0006]   The frequently used FFT has the disadvantage that, due to the block-wise processing of the signals in the time domain, a compromise has to be found between the resolution in the time domain and the resolution in the frequency domain.

[0007]   The frequency of a frequency line is determined according to Equation 1.

$$f(n) = \frac{Fs}{N} \cdot n \tag{1}$$

[0008]   The frequency spacing of the FFT is constant and is obtained from Equation 2.

$$df = \frac{Fs}{N} \tag{2}$$

[0009]   For Fs = 8 kHz and N = 256,

$$df = \frac{8\,kHz}{256} = 31.25\,Hz$$

df      frequency spacing
f       frequency
n       number of the frequency line
Fs      sampling frequency
N       number of frequency lines

[0010]   With a shorter block, for example N = 128, although a better time resolution is obtained, a poorer resolution

is nevertheless obtained in the frequency domain with $df = 62.5\ Hz$. The linear frequency resolution of the FFT thus does not take account of essential psychoacoustic characteristics. By contrast, the frequency resolution of the human ear is nonlinear. The transmission function is described more fully in Eberhard Zwicker: Phychoakustik, Springer Verlag, Berlin, Heidelberg, New York, 1982, pages 20-30. The time resolution of the human ear is approximately 1.9 ms, but that of a 256 point FFT, for example, is 32 ms. Due to these differences between the FFT and the psychoacoustic requirements, a natural-effect speech transmission can be achieved only with limitations in respect of quality. In addition, the additional signal delay due to the block-wise signal processing impairs a telecommunication device both by disrupting the natural flow of a conversation and through the increased echo perception.

[0011] The practice of using a Wiener filter for determining the noise components of a noisy speech signal is generally known. A Wiener filter is described in, for example, "Numerical Recipes in G: The Art of Scientific Computing"; chapter 13.3, Optimal (Wiener) Filtering with the FFT; pages 547-549, Cambridge University Press 1988-1992. With the Wiener filter, the magnitude of the transmission function $|H(n)|$ is calculated for each frequency n, according to Equation 3.

$$|H(n)| = \begin{cases} 1 - o \bullet \left( \dfrac{E(n)}{|X(n)|} \right)^2 & \text{if } |X(n)| > E(n) \\[2em] \text{NFL} & \text{otherwise} \end{cases} \tag{3}$$

$|H(n)|$     magnitude of the transmission function for the frequency n
$E(n)$     estimated averaged value for the ambient noise
$|X(n)|$     magnitude of the noisy speech
NFL     background noise, noise floor
$o$     overestimation factor

[0012] The mean value of the noise is calculated using a first-order recursive filter during the speech pauses. The filter coefficients used are constant.

[0013] According to Equation 3, $|H(n)| = 1$ if $E(n) = 0$, i.e., when there is no noise. If $E(n) \neq 0$, so that the difference becomes less than 1, then, in the ideal case, the noise is subtracted from the spectrum of the noisy speech signal without affecting the speech signal. If, for a frequency n, the power density of the estimated noise $E(n)$ becomes greater than the power density of the estimated noisy speech signal, the above relationship in Equation 3 would produce a negative value. In this case, $|H(n)|$ is set = NFL, so that a background noise NFL is permitted in order to prevent an unnatural masking-out of all noises. The overestimation factor o provided for in Equation 3 serves to reduce errors in the estimation of the energy contents.

[0014] Due to the block-wise processing of the signals by means of the FFT, in the inverse transformation using the IFFT one value is obtained per block, so that a discontinuous value sequence can result which is audible as so-called "musical tones" in the reconverted speech signal. In order to prevent this effect, a sufficiently large value of the background noise NFL is selected to mask the "musical tones". This, however, has the result that only a very limited noise reduction, of approximately 6 dB, can be achieved with the described algorithm and, particularly in the case of a very small speech-to-noise ratio, an improvement is not possible, for example, greater than 10 dB.

[0015] There thus ensues, from the described disadvantages of the noise reduction method using a Wiener filter, the object of altering the noise estimation by means of the Wiener filter and the rules for transforming the noisy speech signals from the time domain into the frequency domain and vice versa so as to permit an adaptation to the nonlinear transmission behaviour of the human ear.

[0016] This object is achieved by the method disclosed in the first claim.

[0017] The essence of the invention consists in that the conditions for determining the transmission function of the Wiener filter are optimized and that a Continuous Fourier Transformation is used as a rule for transforming the noisy speech signal. The Continuous Fourier Transformation is described in the patent application DE 10 111 249.1.

[0018] The application of the Continuous Fourier Transformation creates new conditions for an improved noise reduction.

[0019] The application of the rule, described in connection with Equation 3, for the transmission function $|H(n)|$ of the Wiener filter of the prior art has the result that, in the case of small speech signals, $|H(n)|$ becomes = NFL and,

consequently, speech syllables with a low energy content are omitted from the output signal. The sum of the speech signal and noise |X(n)| is a highly modulated signal which exceeds the noise level E(n) only temporarily, when the energy of the corresponding frequency of the speech signal is just in the transition to the energy content of the noise threshold value. This effect occurs particularly when the noise is modulated and superimposed on the speech signal.

[0020]   In order to achieve a greater sensitivity for small speech signal-to-noise ratios, the changeover of the transmission function |H(n)| to the background noise NFL is only permitted, according to the invention, if the estimated mean value of the speech signal SE(n) is not greater than the estimated mean value of the noise E(n), see Equation 4.

$$|H(n)| = \begin{cases} 1 - o \cdot \left( \dfrac{E(n)}{|X(n)|} \right)^2 & \text{if } |X(n)| > E(n) \\[20pt] \text{NFL} & \text{if } \overline{SE(n)} > \overline{E(n)} \end{cases} \qquad (4)$$

[0021]   Due to this rule, even faint components of the speech signal are reliably transmitted, and the system is thus better adapted to the speech spectrum.

[0022]   A first-order recursive filter permits determination of the estimated mean values of the Speech signal SE(n) and of the noise E(n). The speech signal SE(n) is estimated during the speech activity, pause indicator p = 0, and the noise E(n) is estimated during the speech pauses, pause indicator p = 1, according to Equations 5 and 6.

$$SE(n,k) = \begin{cases} \alpha(n) \cdot |X(n,k)| + \beta(n) \cdot SE(n,k-1) & \text{if } p = 0 \\[16pt] SE(n,k-1) & \text{otherwise} \end{cases} \qquad (5)$$

$$E(n,k) = \begin{cases} \alpha(n) \cdot |X(n,k)| + \beta(n) \cdot E(n,k-1) & \text{if } p = 1 \\[16pt] E(n,k-1) & \text{otherwise} \end{cases} \qquad (6)$$

k        sampling instant
p        pause indicator
$\alpha, \beta$     filter coefficients, which can assume fixed values or be frequency-dependent

[0023]   The values SE(n) and E(n) determined according to Equations 5 and 6 are calculated in dependence on frequency and produce an optimum time response.

[0024]·  In order to prevent disturbing transient noise fluctuations, Equation 3 is expanded in such a way that the difference is only formed if the speech signal SE(n) is greater than the noise E(n), see Equation 4. The time response of the speech signal SE(n) can then be determined according to the speech characteristics, which differ from short excitations of the noise E(n).

$$|H(n)| = \begin{cases} 1 - o \cdot \left(\dfrac{E(n)}{|X(n)|}\right)^2 & \text{if } (SE(n) > E(n)) \,\&\, (|X(n)| > E(n)) \qquad (7) \\ \\ NFL & \text{if } \overline{SE(n)} > \overline{E(n)} \end{cases}$$

[0025] The unwanted "musical tones" effect of the known noise reduction methods is eliminated if, instead of the transformation methods such as, for example FFT and IFFT, which work in blocks, transformation methods are used in which the nonlinear frequency resolution of the human ear is taken into account. Thus, a range of auditory characteristics, such as frequency resolution, time resolution and selection characteristics must be taken into account if a natural-sounding speech signal, or an audio signal generally, is to be received. In order to achieve this, a Fourier transformation has already been disclosed which is adapted to the transmission function of human sensory organs, cf. DE 101 11 249.1. This transformation deviates from the fixed assignment of number of frequencies N equal to number of sampling values K, which necessitate a constant frequency spacing according to Equation 1 and a constant bandwidth B, and a Continuous Fourier Transformation CFT and an Inverse Continuous Fourier Transformation ICFT of the speech are performed. In the case of the CFT, a time function x(k) is mapped in frequency groups, the number and magnitude of which are determined, for example, according to the BARK scale, cf. Kapust, Rolf: Qualitätsbeurteilung codierter Audiosignale mittels einer BARK-Transformation, Dissertation 1993, University of Erlangen-Nümberg. Within a frequency group, a number of frequency lines N is calculated so that the frequency resolution and the time resolution are matched to the transmission function of the human ear. The bandwidth B(n) with which a frequency line is transmitted is determined from the frequency lines n+1 and n-1 adjacent to a frequency line n. From the bandwidth B(n) is determined the limiting frequency fg of a low-pass filter which, as an integrator, replaces the otherwise usual summation of the blocks and thus effects a sliding transformation. A rapid modification and, consequently, an adaptation to the current situation of the calculated transmission function |H(n)| is already achieved with 17 frequency lines, at a sampling rate of 8 kHz. This rapid modification results in a modulation of the reconverted speech. An improved time response of the transmission function |H(n)| is achieved if a frequency-dependent short average magnitude SAM (|H(n)|) of the transmission function is formed, and a noise-reduced frequency line n is thus produced. The short average magnitude SAM (|H(n)|) is formed using a recursive filter such as that described in, for example, EP 1 005 016 A2 and represented in Fig. 3 thereof.

[0026] The low-pass used as an integrator in the case of the Continuous Fourier Transformation CFT for the purpose of determining each frequency line can be further improved in the formation of the complex frequency, for the purpose of improving the speech quality in noise reduction systems. Since speech signals exist for a certain duration, for example, longer than 100 ms, and noises can nevertheless occur in shorter time intervals during the speech, it is useful to determine a real component and an imaginary component of the complex frequency according to Equations 8, 9 and 10. Equations 8 and 9 describe a first-order recursive low-pass filter.

$$re(n,k) = \cos(n,k) \cdot x(k) \cdot \alpha x(n) + re(n,k-1) \cdot \beta x(n) \qquad (8)$$

$$im(n,k) = \sin(n,k) \cdot x(k) \cdot \alpha x(n) + im(n,k-1) \cdot \beta x(n) \qquad (9)$$

the filter coefficients x(n) being determined according to the following Equation 10.

$$x(n) = \begin{cases} \kappa \cdot \tau & \text{if } |re(k) - im(k)| > |re(k-1) - im(k-1)| \\ \tau & \text{otherwise} \end{cases}$$

$$\tau = \frac{1}{2 \cdot \pi \cdot fb(n)} \qquad fb = \text{bandwidth of the frequency line} \qquad (10)$$

$$\kappa = 2.....10 \quad const.$$

[0027] This modification has the effect that interruptions in the speech signal due to reduction of very large, short noises are restored. Due to the large time constant effected by the filter coefficient x(n), the current magnitude and the current phase position are maintained, so that speech interruptions are avoided.

[0028] If a large noise reduction is to be achieved, the background noise NFL assumes a very small value. This also results in the suppression of very weak speech signals, which may then be evaluated as noise. In order to prevent this effect, the background noise can be determined in dependence on the current requirements, according to Equation 11.

$$nfl(n, k) = \begin{cases} nava(n) \bullet NFL + navb(n) \bullet nfl(n, k-1) & \text{if } SE(n) > E(n) \\ NFL_{min} & \text{otherwise} \end{cases} \qquad (11)$$

$$nava(n) = 1 - navb(n) = 1 - e^{\frac{2 \cdot \pi \cdot fb(n)}{Fs}}$$

Fs = sampling frequency
Fb(n) = bandwidth of the frequency line n

nava noise floor average a
navb noise floor average b

[0029] Equation 11 is used to average a background noise nfl(n), which is dependent on the frequency, if the speech signal SE(n) is greater than the noise E(n). When speech is present, the value for the background nfl(n) is greater than the minimum background noise, so as to ensure that speech signals are not suppressed.

[0030] The overestimation factor o determines the magnitude of the noise reduction during the speech activity. A large noise reduction requires a small overestimation factor o. Experiments have shown that an optimum overestimation factor o can be determined according to Equation 12.

$$o(n) = \frac{1}{\log(nfl(n))} \qquad (12)$$

[0031] Taking into account the conditions, adapted to the nonlinear transmission behaviour of the human ear, for determining the transmission function ($|H(n)|$) of the Wiener filter, then

$$|H(n)| = SAM(n) \begin{cases} 1 = o(n) \bullet \left(\frac{E(n)}{|X(n)|}\right)^2 & \text{if } (SE(n) > E(n)) \& (|X(n)| > E(n)) \\ nfl(n) & \text{if } (SE(n) < E(n)) \end{cases} \qquad (13)$$

[0032] With this rule, the nonlinear transmission behaviour of the human ear is taken into account. Despite a large noise reduction, loss of quality in the speech transmission is prevented by means of the method.

[0033] The invention is explained further with reference to an embodiment example and the associated drawing,

wherein:

Fig. 1     shows a block diagram of a circuit arrangement for spectral subtraction using a Wiener filter according to the prior art,

Fig. 2     shows a block diagram of a circuit arrangement for spectral subtraction using a Wiener filter and application of a Continuous Fourier Transformation,

Fig. 3     shows a block diagram for the application of the Continuous Fourier Transformation for the purpose of reducing noise, and

Fig. 4     shows a distribution of the frequency lines to the frequency groups in the case of the Continuous Fourier Transformation.

[0034] As shown by Figure 1, a circuit arrangement for noise reduction consists essentially of two modules for windowing 1.1, 2.1 of the analog-digital converted input signal x(k), a speech detector 1.2, two noise averaging devices 1.3, 2.3, two Wiener filters 1.4, 2.4 and an overlap add 1.5, as well as the modules for the Fast Fourier Transformation FFT 1.6, 2.6 and for the Inverse Fast Fourier Transformation 1.7, 2.7. For the purpose of processing the input signal x(k) by means of the FFT, the input signal x(k) is divided into blocks, of the length N, also called windows, in such a way that the spectral characteristics are largely constant for the duration of the window. Whereas, in the middle of the window, the course of the function can be precisely described, the information on how the function continues is absent at the edge of the window. Two windows, offset by $\frac{1}{2}N$, are therefore processed, for example, according to the Hamming function and, following back-transformation, overlapped by means of an overlap add 1.5 so that the energy values are not falsified at the edges of the windows. The noise averaging device 1.3, 2.3 is used to determine a mean value, in the speech pauses, from the input signal x(k) transformed into the frequency domain. The speech pause is ascertained by a speech detector 1.2 which delivers a signal p as a pause indicator, p = 0 corresponding to speech, p = 1 corresponding to speech pause. The power density of the noise spectrum H(n) is calculated using the Wiener filter 1.4, 2.4 and subtracted from the noisy speech signal X(n), so that the noise-corrected speech signal SE(n) can be transformed back out of the frequency domain into the time domain by means of the IFFT and, following overlapping of the windows, the speech signal y(k) is formed in the time domain.

[0035] The disassociation from block processing in the FFT and IFFT renders windowing and window overlapping superfluous, as shown in Fig. 2. Otherwise, the method steps described in connection with Fig. 1 are also performed in the application of the Continuous Fourier Transformation CFT and the Inverse Continuous Fourier Transformation ICFT according to Fig. 2.

[0036] Fig. 3 shows an example for the application of the CFT/ICFT. The input signal x(k) is divided into four frequency groups, scaled logarithmically. This division is effected, for example, at a sampling frequency Fs = 8 kHz, there being formed a first frequency group with a bandwidth B = 500 Hz, at a first sampling frequency $\frac{1}{8}Fs = 1000Hz$, a second frequency group with a bandwidth B = 1000 Hz, at a second sampling frequency $\frac{1}{4}Fs = 2000Hz$, a third frequency group with a bandwidth B = 2000 Hz, at a third sampling frequency $\frac{1}{2}Fs = 4000Hz$, and a fourth frequency group for frequencies over 2000 Hz, at the sampling frequency Fs = 8 kHz. Via the bandpass filters BP 500, BP 1000 and BP 2000, and via the high-pass filter HP 2000, the input signal x(k) according to Fig. 3 is transformed by means of the CFT into the frequency domain, in which it is processed according to the application and transformed back into the time domain, as y(k), by means of the ICFT, via low-pass filters LP and interpolation filters IP and through summation of the frequency groups.

[0037] Fig. 4 shows the distribution of the frequency lines to the frequency groups, as is particularly advantageous, for example, in the case of an economically optimized version. This distribution is eminently suitable in the case of the application of noise reduction in the spectral domain. The first frequency group up to 500 Hz is allotted 40 frequency lines, the second frequency group up to 1000 Hz is allotted 20 frequency lines, the third frequency group up to 2000 Hz is allotted 10 frequency lines and the fourth frequency group up to 4000 Hz is allotted 5 frequency lines. In the noise reduction example illustrated, a high frequency resolution is desired in precisely that frequency range in which the majority of frequencies which are attributable to the interfering noise occur, i.e., practically, the range between f = 0 and 2 kHz. As shown in Fig. 4, 75 frequency lines have been logarithmically distributed such that the frequency resolution in the lower frequency range up to 500 Hz is particularly high, in this case being 10 Hz. Such a frequency resolution is not even achieved with a FFT with 512 frequency lines, the frequency resolution in this case being 16 Hz. As shown by Fig. 4, the frequency resolution decreases, to the topmost frequency line, to 510 Hz, corresponding to a time resolution of 0.98 ms, whereas the FFT with 512 frequency lines has a constant value of 31.25 ms. The necessary computational requirement can be greatly reduced through subsampling with decimation filters and interpolation filters. The range with the most frequency lines can be subjected to the greatest subsampling. Experiments have shown that

the above-mentioned 75 frequency lines per sampling value can be reduced to 20 frequency lines per sampling value without loss of quality of a natural-sounding speech.

## Claims

1. Method for improving noise reduction in speech transmission by applying a rule for transforming a noisy speech signal in the time domain into a noisy signal in the frequency domain and using a Wiener filter with the transmission function

$$
\left|H(n)\right| = \begin{cases} 1 - o \bullet \left( \dfrac{E(n)}{\left|X(n)\right|} \right)^2 & \text{if } (\left|X(n)\right| > E(n)) \\[4mm] \text{NFL} & \text{otherwise} \end{cases}
$$

for evaluating the noise spectrum for the purpose of performing a spectral subtraction of the noise spectrum (E(n)) from the frequency spectrum of the noisy speech signal,
**characterized in that,**
for the transmission function H(n), the value of a background noise NFL is set if the estimated mean value of the speech signal (SE(n)) is smaller than the estimated mean value of the noise (E(n)),
**in that** for the transmission function H(n), a current value is calculated for a frequency if the mean value of the speech signal (SE(n) is greater than the estimated mean value of the noise (E(n)) and the magnitude of the noisy speech signal |X(n)| is greater than the estimated mean value of the noise (E(n))
and **in that**, in application of a Continuous Fourier Transformation for the transformation of the noisy speech signal from the time domain into the frequency domain, a frequency-dependent short average magnitude (SAM(n) is formed for the transmission function H(n).

2. Method according to Claim 1, **characterized in that** the value of the background noise is calculated for a frequency in dependence on the noise reduction factor and in dependence on the probability with which this frequency occurs in the speech spectrum.

3. Method according to Claim 1, **characterized in that** the value of an overestimation factor o is selected which is equal to the reciprocal value of the decimal logarithm from the noise reduction factor.
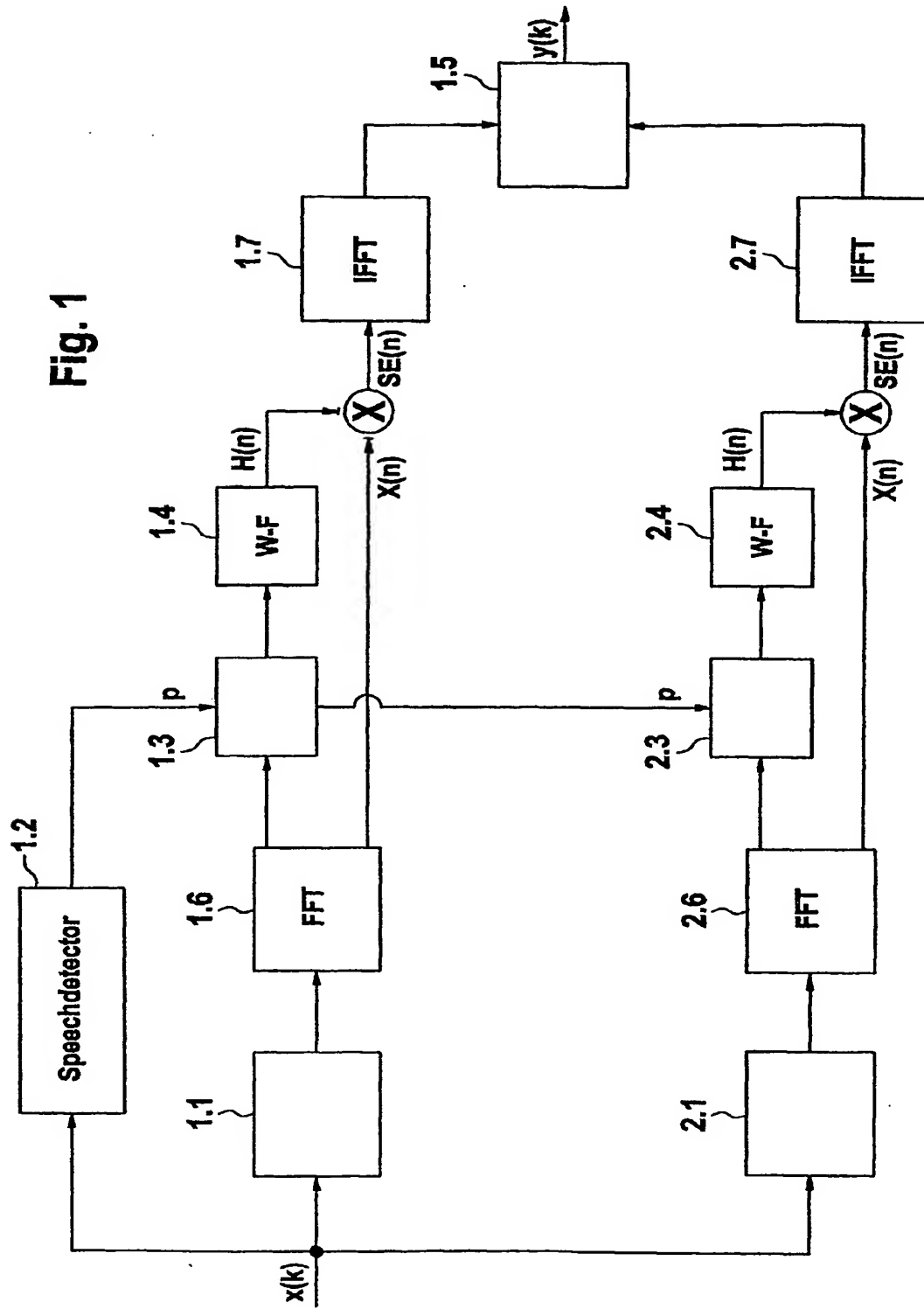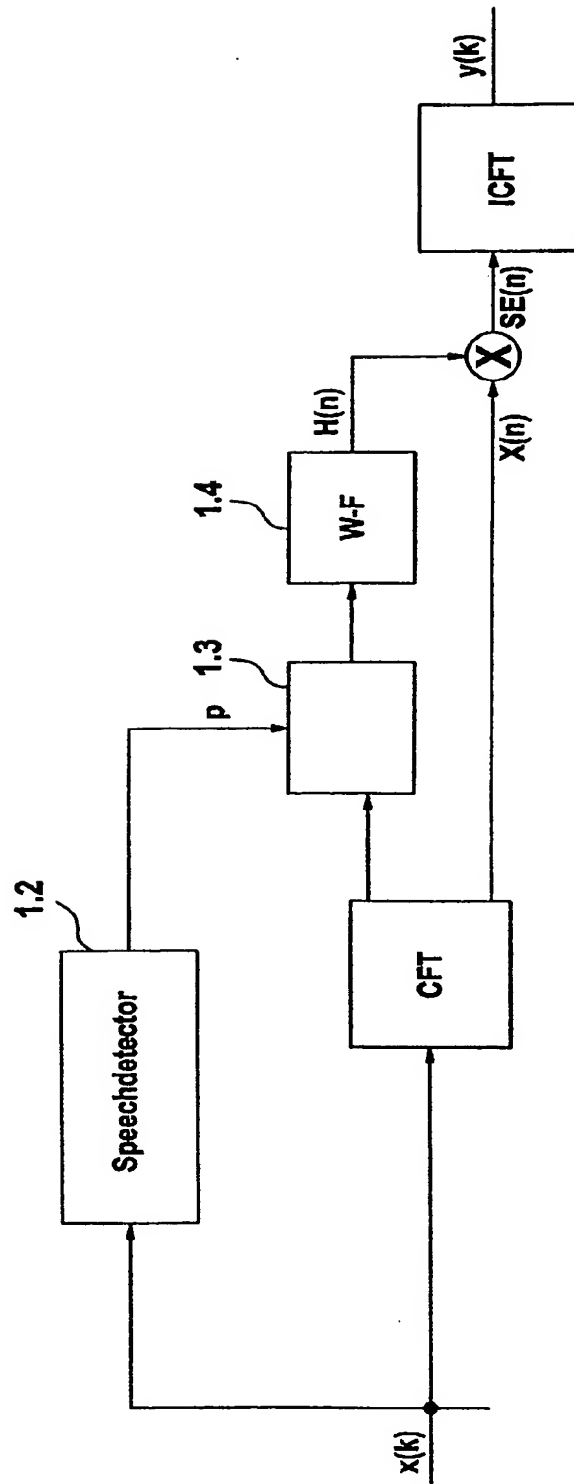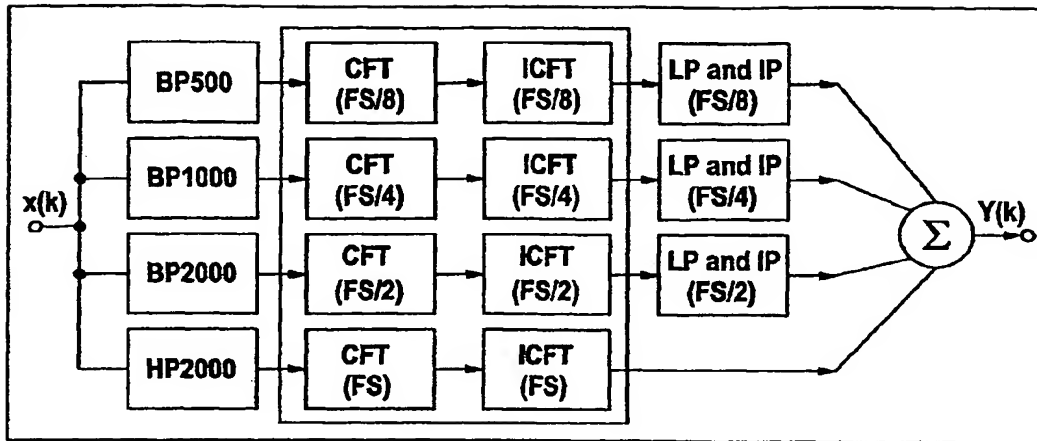
## Fig. 1

## Fig. 2

# Fig. 3



# Fig. 4